

HIGH-PERFORMANCE PARALLEL INTERFACE - 6400 Mbit/s Physical Switch Control (HIPPI-6400-SC)

February 28, 1997

Secretariat:

Information Technology Industry Council (ITI)

ABSTRACT: HIPPI-6400-SC provides a protocol for controlling physical layer switches which are based on the High-Performance Parallel Interface at 6400 Mbits/s (HIPPI-6400-PH), a simple high-performance point-to-point interface for transmitting digital data at peak data rates of 6400 Mbit/s between data-processing equipment.

NOTE:

This is an internal working document of X3T11, a Technical Committee of Accredited Standards Committee X3. As such, this is not a completed standard. The contents are actively being modified by X3T11. This document is made available for review and comment only. For current information on the status of this document contact the individuals shown below:

POINTS OF CONTACT:

Roger Cummings (X3T11 Chairman)
Distributed Processing Technology
140 Candace Drive
Maitland, FL 32751
(407) 830-5522 x348, Fax: (407) 260-5366
E-mail: cummings_roger@dpt.com

Ed Grivna (X3T11 Vice-Chairman)
Cypress Semiconductor
2401 East 86th Street
Bloomington, MN 55425
(612) 851-5200, Fax: (612) 851-5087
E-mail: elg@cypress.com

Roger Ronald (HIPPI-6400-SC Technical Editor)
E-Systems
MS 35300 HD
PO Box 660023
Dallas, TX 75266-0023
(214) 205-8043, Fax: (214) 272-8144
E-mail: rronald@esy.com

Comments on Rev 0.20

This is a preliminary document. The first draft (rev 0.01) was presented and reviewed for the first time in March 1996. The second revision (rev .10) was reviewed on May 9th and 10th in Dallas. This revision corrects errors discovered at that time and continues the process of documentation.

Rev bars are now included in this revision of the document except for cases of minor punctuation or spelling error correction.

Major changes from the previous revision include:

- new definitions added for alternate pathing, final destination and original source
- changing the requirements for selection of the switch port to allow alternate pathing
- beefing up the paragraph on error checking required in fabric
- adding information on micropacket interleaving
- clarifying the paragraph on congestion management
- adding reserved addresses from the HIPPI-800-SC standard
- adding more information on routing to the informative appendix

In addition to discussing these changes, it is expected that the next meeting (Santa Fe, June 10-11) will cover usage of administrative packets.

Comments on Rev 0.30

This revision was started to collect changes and additions made during and after the June ANSI meeting held in Santa Fe, NM on June 10th thru the 12th.

Major changes from the previous revision include:

- Added definitions
- General clean-up
- Pruning of sections detailing alternative addressing approaches (alternate pathing, broadcast, and multicast)
- Moved congestion management paragraph to reside within the general section on error protection
- Removed requirement of 4 micropacket message support for in-band communications
- Split switching, bridging, and routing into three appendices while adding text and examples

Comments on Rev 0.40

This revision was started to collect changes and additions made during and after the July HIPPI-6400 working group meeting held in San Jose on July 11th and 12th.

Major changes from the previous revision include:

- Definitions added for administrator, fabric, log, and switch
- Switch addressing references to optional modes reduced to a minimum
- Address restrictions for inter-operation with HIPPI-800 removed from section 6
- Removed e-mail list instructions from this page for obsolete mail groups

Comments on Rev 0.45

This revision was started to collect changes and additions made during and after the August HIPPI-6400 ANSI meeting held in Honolulu on August 5-7, 1996. Because the document has not been reviewed line-by-line since the July working group meeting, change bars still include the revision 0.4 changes.

Major changes from the previous revision include:

- Updated definitions and acronyms to follow the lead of HIPPI-6400-PH
- Removal of comments that this specification would describe switch-to-switch negotiation of address configuration.
- Information and procedures for using admin micropackets for topology discovery.
- Information and procedures for using admin micropackets for logical address assignment.
- Replacement of 16 bit logical addresses with 48 bit Universal LAN Addresses (ULAs) and provision for optional operation using 16 of the 48 bits.
- Removed requirement to support 64K switch addresses.
- Changed the limit for the maximum count of micropackets that may be sent on a single VC before interleaving traffic from other VCs from 65 to 66 (to match the limit for a VC0 message).
- Updated text to reflect decision that all micropackets except Header micropackets (not just Data micropackets) will be treated as part of a message following a Header micropacket.

Comments on Rev 0.50

This revision was started to collect changes and additions made during and after the September and October HIPPI-6400 ANSI meetings. Because the document has not been reviewed line-by-line since the July working group meeting, change bars still include the revision 0.4 changes.

There are no major changes from the previous revision.

Comments on Rev 0.60

This revision was started to collect changes and additions made during and after the November HIPPI-6400 working group meeting held in Phoenix on November 6th and 7th.

Major changes from the previous revision include:

- Inclusion of the Admin Micropacket Draft, revision 0.4. This inclusion does not add change bars as that draft had been reviewed by the group.
- Changed admin command and response names to include an underline between multiple words. This change does NOT have change bars.
- Added many “shalls” to the admin micropacket command and response table.
- Added a diagram showing address processing

Plus, this version continued to clean up document editorial comments

Comments on Rev 0.70

This revision was started to collect changes and additions made during and after the December HIPPI-6400 meeting held in Minneapolis on December 2nd and 3rd.

Major changes from the previous revision include:

- Changing the undefined admin micropacket Function codes to reserved.
- Providing a admin micropacket command/response for endpoints to learn connected addresses in support of broadcast capability (ULA_LIST_REQUEST and ULA_LIST_RESPONSE).
- Changed the name of the RETURN_LOGICAL_ADDRESS and LOGICAL_ADDRESS_RESPONSE to ULA_REQUEST and ULA_RESPONSE.
- Moved the text on admin micropacket commands and status out of the overly large table.
- Deleted Annex C on routing (no useful information was included in that section).

Comments on Rev 0.80

This revision was started to collect changes and additions made during and after the January HIPPI-6400 meeting held in Phoenix on January 7th, 8th and 9th.

Major changes from the previous revision include:

- The rules for address matching of administrative micropackets in section 7.6 and Figure 5 were modified to prevent an element with an address of x'FFFFFFFF' from taking and processing all micropackets sent using hop-count addressing.
- Change section 8 title from "Address Configuration" to "ULA Configuration"
- Added a bibliography as Annex C

Comments on Rev 0.90

This revision was started to collect changes and additions made during and after the February HIPPI-6400 meeting held in San Jose.

Major changes from the previous revision include:

- Capitalized "Admin" (no change bars added)
- The document was scrubbed to eliminate use of the word "address" by itself. All occurrences were changed to be "element addresses" or ULAs.

For broadcast functionality, the following was added:

- The EXCHANGE_TYPE and TYPE_RESPONSE now have two bits to indicate whether a endpoint desires to receive broadcast messages and whether the endpoint is willing to be a broadcast server. EXCHANGE_TYPE is now required to be sent at one second intervals for endpoints who want broadcasts and/or are willing to be servers. This serves as a ping function to allow updating of who the broadcast server is and who should receive broadcasts.
- Specified that the first ULA registered on a port using the ULA_REQUEST/ULA_RESPONSE is the destination ULA that will be used for broadcast messages.
- The ULA_RESPONSE now contains a broadcast address that may be used by hosts.
- Made clear that exactly one ULA is returned for each port in the ULA_LIST_RESPONSE.
- Added that switches must do type discovery and ULA exchange.
- Changed 8.2.2 to describe switch to switch ULA operations.
- Added clause 9 to describe how broadcast works.

Please help us in this development process by sending comments, corrections, and suggestions to the Technical Editor, Roger Ronald @ E-Systems via e-mail (rronald@esy.com).

Table Of Contents

1	Scope	1
2	Normative references	1
3	Definitions and conventions	1
3.1	Definitions	1
3.2	Editorial conventions	2
3.2.1	Binary notation	2
3.2.2	Hexadecimal notation	2
3.2.3	Acronyms and other abbreviations	2
4	System overview	3
4.1	Switch function	3
4.2	Micropacket	3
4.3	Message	4
5	Switch routing	4
5.1	Micropacket data transferred through fabric	4
5.2	Routing of Header micropacket	4
5.2.1	Switch addressing	4
5.3	Routing of subsequent micropackets in a Message	5
5.4	Error protection	5
5.4.1	Mandatory error checking	5
5.4.2	Optional error checking	5
5.4.3	Congestion management	5
5.5	Data interleaving	5
5.5.1	Micropacket interleaving	5
5.5.2	Message interleaving	5
6	ULA restrictions and reserved ULAs	6
7	Admin micropackets	7
7.1	Element definition	8
7.2	Element conformance	8
7.3	Admin micropacket functions	8
7.4	Admin micropacket format	9
7.5	Admin micropacket functions	10
7.5.1	PING	10
7.5.2	PING RESPONSE	10
7.5.3	SET_ELEMENT_ADDRESS	10
7.5.4	SET_ELEMENT_ADDRESS_RESPONSE	12
7.5.5	RESET	12
7.5.6	EXCHANGE_ELEMENT_FUNCTION	12
7.5.7	ELEMENT_FUNCTION_RESPONSE	12
7.5.8	ULA_REQUEST	13
7.5.9	ULA_RESPONSE	13
7.5.10	READ_REGISTER	13
7.5.11	READ_REGISTER_RESPONSE	13
7.5.12	WRITE_REGISTER	13
7.5.13	WRITE_REGISTER_RESPONSE	14
7.5.14	ERROR_RESPONSE	14
7.5.15	ULA_LIST_REQUEST	14
7.5.16	ULA_LIST_RESPONSE	14
7.5.17	Reserved Admin micropacket functions	14
7.6	Addressing of Admin micropackets	14
7.7	Admin Element address assignment	15
7.8	Admin micropacket flow control	15

8	ULA Configuration	16
8.1	Determination of Topology	16
8.2	ULA exchange	17
8.2.1	Endpoints on both ends	17
8.2.2	Switches on both ends	17
8.2.3	Endpoint to switch	18
9	Broadcast	18
9.1	Broadcast Operation	18
9.2	Selection of broadcast server.	18
9.3	Registration for broadcast	19
9.3.1	Endpoints	19
9.3.2	Switches	19
A	Switching	20
A.1	General	20
A.2	Logical addressing	20
A.3	Input specific logical addressing	21
B	Bridging	22
B.1	General	22
C	Bibliography	23

List Of Figures

Figure 1 -	Message format	4
Figure 2 -	Header micropacket addressing.....	4
Figure 3 -	HIPPI-6400 Switch	6
Figure 4 -	Admin Micropacket Byte Format	8
Figure 5 -	Admin micropacket addressing	9
Figure 6 -	Endpoint to endpoint connect.....	17
Figure 7 -	Hosts and switch configuration.....	20
Figure 8 -	Hosts, switch, and bridge configuration.....	22

List Of Tables

Table 1 -	Data carried through fabric.....	3
Table 2 -	Data to route 1st micropacket in a Message	3
Table 3 -	Data to Route Subsequent Micropackets in a Message	3
Table 4 -	Data used for error checking and reporting.....	3
Table 5 -	Reserved ULAs	6
Table 6 -	Admin Micropacket Format	9
Table 8 -	Status Flags	10
Table 7 -	Admin Commands.....	11
Table 9 -	Port look-up table	21

Foreword (This Foreword is not part of American National Standard X3.xxx-199x.)

This American National Standard specifies the behavior and control for HIPPI-6400 physical layer switches. HIPPI-6400 is an efficient high-performance point-to-point interface. HIPPI-6400 physical layer switches may be used to give the equivalent of multi-drop capability, connecting together multiple data processing equipments.

This standard provides an upward growth path for legacy HIPPI-based systems.

This document includes annexes which are informative and are not considered part of the standard.

Requests for interpretation, suggestions for improvement or addenda, or defect reports are welcome. They should be sent to the X3 Secretariat, Information Technology Industry Council, 1250 Eye Street, NW, Suite 200, Washington, DC 20005.

This standard was processed and approved for submittal to ANSI by Accredited Standards Committee on Information Processing Systems, X3. Committee approval of the standard does not necessarily imply that all committee members voted for approval. At the time it approved this standard, the X3 Committee had the following members:

(List of X3 Committee members to be included in the published standard by the ANSI Editor.)

Subcommittee X3T11 on Device Level Interfaces, which developed this standard, had the following participants:

(List of X3T11 Committee members, and other active participants, at the time the document is forwarded for public review, will be included by the Technical Editor.)

Introduction

This 6400 Mbits/second High-Performance Parallel Interface, Physical Switch Control (HIPPI-6400-SC) standard defines the control for HIPPI-6400 physical layer switches. HIPPI-6400 is an efficient high-performance point-to-point interface. Small fixed-size micropackets provide an efficient, low-latency, structure for small messages, and a building block for large messages. HIPPI-6400 physical layer switches may be used to give the equivalent of multi-drop capability, connecting together multiple data processing equipments.

Characteristics of this HIPPI-6400 physical switch control protocol include:

- Support for 48 bit Universal LAN Addresses (ULAs)
- Support for restricted mode operation with a 16 bit subset of the ULA
- Procedures for use of Admin micropackets to automate ULA assignment
- Ability to span multiple physical layer switches within a fabric
- Support for physical layer switches with differing numbers of ports, all within the same fabric
- Specified reserved ULAs to aid address self-discovery, switch management, and switch control
- Support for 4 Virtual Channels
- Support for broadcast capabilities, either within a switch or provided by an attached server

American National Standard for Information Technology –

High-Performance Parallel Interface – 6400 Mbit/s Physical Switch Control (HIPPI-6400-SC)

1 Scope

This American National Standard provides switch control for physical layer switches using the 6400 Mbits/second High-Performance Parallel Interface (HIPPI-6400), a high-performance point-to-point interface between data-processing equipment.

The purpose of this standard is to facilitate the development and use of the HIPPI-6400 in computer systems by providing common physical switch control. The standard provides switch control structures for physical layer switches interconnecting computers, high-performance display systems, and high-performance, intelligent block-transfer peripherals. This standard also applies to point-to-point HIPPI-6400 topologies.

Specifications are included for:

- Interleaving of Virtual Channels (VCs) within a physical channel
- Selection of Messages for transmission on physical channels
- Self discovery of configuration information

2 Normative references

The following American National Standard contains provisions which, through reference in this text, constitute provisions of this American National Standard. At the time of publication, the edition indicated was valid. All standards are subject to revision, and parties to agreements based on this standard are encouraged to investigate the possibility of applying the most recent edition of the standard listed below.

ANSI X3.183-1991, *High-Performance Parallel Interface – Mechanical, Electrical, and Signalling Protocol Specification (HIPPI-PH)*.

ANSI X3.210-1992, *High-Performance Parallel interface, Framing Protocol (HIPPI-FP)*.

ANSI X3.222-1993, *High-Performance Parallel interface, Physical Switch Control (HIPPI-SC)*.

ANSI X3.xxx-199x, *High Performance Parallel Interface 6400 Mbits/s, Physical Layer (HIPPI-6400-PH)*

3 Definitions and conventions

3.1 Definitions

For the purposes of this standard, the following definitions apply.

3.1.1 Admin micropacket: A special HIPPI-6400 micropacket used for station management.

3.1.2 administrator: A station management entity providing external management control.

3.1.3 alternate pathing: Capability to address a Message to select from a group of ports based upon defined criteria.

3.1.4 broadcast: The capability for a Source to send one message that arrives at multiple Destinations.

3.1.5 Destination: The equipment that receives the data.

3.1.6 Device: Any system level component (e.g. endpoint or switch) with a HIPPI-6400 port.

3.1.7 Element: Any distinct part of a HIPPI-6400 system that is able to receive and send Admin micropackets conforming to this standard.

3.1.8 Element address: A 32 bit field specifying where an Admin micropacket originated or is delivered.

3.1.9 fabric: All of the switching equipment connected together in a configuration.

3.1.10 Final Destination: The end device that receives, and operates on, the payload portion of the micropackets. This is typically a host computer system, but may also be a translator, bridge, or router.

3.1.11 HIPPI-PH: High-Performance Parallel Interface - Mechanical, Electrical, and Signalling Protocol Specification (HIPPI-PH), ANSI X3.183-1991. Data is transmitted in parallel over copper twisted-pair cables at 800 or 1600 Mbits per second.

3.1.12 HIPPI port: A HIPPI-6400-PH, or HIPPI-PH, Source or Destination.

3.1.13 in-band: Switch control communications accomplished over a HIPPI-6400 link. As opposed to out-of-band (using an alternative communication channel).

3.1.14 link: A full-duplex connection between HIPPI-6400-PH Devices.

3.1.15 log: The act of making a record of an event for later use.

3.1.16 micropacket: The basic transfer unit consisting of 32 data bytes and 64 bits of control information.

3.1.17 Message: An ordered sequence of one or more micropackets which have the same VC. The first micropacket is a Header micropacket. The last micropacket, which may also be the first micropacket, has the TAIL bit set.

3.1.18 optional: Characteristics that are not required by HIPPI-6400-SC. However, if any optional characteristic is implemented, it shall be implemented as defined in HIPPI-6400-SC.

3.1.19 Originating Source: The end device that generates the payload portion of the micropackets. This is typically a host computer system, but may also be a translator, bridge, or router.

3.1.20 Source: The equipment that transmits the data.

3.1.21 switch: An equipment that provides connections between HIPPI-6400 links based on this standard.

3.1.22 Universal LAN Address: A logical address stored in a Source or Destination field that uniquely

identifies an Originating Source or Final Destination. The ULA conforms to the format of other networking protocols (e.g. Ethernet).

3.1.23 Virtual Channel (VC): One of four logical paths within each direction of a link.

3.2 Editorial conventions

In this standard, certain terms that are proper names of signals or similar terms are printed in uppercase to avoid possible confusion with other uses of the same words (e.g., FRAME). Any lowercase uses of these words have the normal technical English meaning.

A number of conditions, sequence parameters, events, states, or similar terms are printed with the first letter of each word in uppercase and the rest lowercase (e.g., State, Source). Any lowercase uses of these words have the normal technical English meaning.

The word *shall* when used in this American National standard, states a mandatory rule or requirement. The word *should* when used in this standard, states a recommendation.

3.2.1 Binary notation

Binary notation is used to represent relatively short fields. For example a two-bit field containing a binary value of 10 is shown in binary format as b'10'.

3.2.2 Hexadecimal notation

Hexadecimal notation is used to represent some fields. For example a two-byte field containing a binary value of b'1100010000000011' is shown in hexadecimal format as x'C403'.

3.2.3 Acronyms and other abbreviations

ACK	acknowledge indication
ARP	Address Resolution Protocol
CR	credit amount parameter
CRC	cyclic redundancy check
ECRC	end-to-end CRC
HIPPI	High-Performance Parallel Interface
IP	Internet Protocol
LCRC	link CRC
MAC	Media Access Control
ns	nanoseconds
RIP	Routing Information Protocol

RSEQ	receive sequence number
TSEQ	transmit sequence number
ULA	universal LAN address
VC	virtual channel
VCR	virtual channel credit selector
μs	microseconds

4 System overview

This paragraph provides an overview of the structure, concepts, and mechanisms used in HIPPI-6400-SC.

4.1 Switch function

HIPPI-6400 switches provide a method to send Messages from a Source port to a Destination port. Each Message travels on one of the four Virtual Channels (VCs) available in HIPPI-6400-PH (see HIPPI-6400-PH for assignments of Message type to VC). All of the micropackets of a Message are transmitted on a single VC, i.e., the VC number does not change as the micropackets travel from the Originating Source to the Final Destination over one or more links.

Different VCs are interleaved on the physical channel allowing up to four Messages to proceed to a Destination or from a Source at any given time.

During transfer of a Message, the VC in use is busy and is unavailable for use by other Messages involving the same Source or Destination ports.

4.2 Micropacket

Micropackets are the basic transfer unit for HIPPI-6400. As described in HIPPI-6400-PH, a micropacket is composed of 32 data bytes and 64 bits of control information.

The 64 bits of control information in each micropacket includes parameters for physical (PH) layer functions and for switch control (SC) functions. These functions include:

- selecting a VC
- detecting missing micropackets
- denoting the types of information in the micropacket

- marking the last micropacket of a Message
- signalling that the Message was truncated at its originator, or damaged en-route, and should be discarded

Table 1 describes the information that the switch fabric carries from a HIPPI-6400-PH source to a HIPPI-6400-PH destination. Table 2 and table 3

Table 1 - Data carried through fabric

Description	Size
ERROR	1 Bit
TAIL	1 Bit
VC	2 Bits
TYPE	4 Bits
ECRC	16 Bits
Payload Data	32 Bytes

describe the information that a switch fabric uses to determine micropacket routing.

Table 2 - Data to route 1st micropacket in a Message

Description	Size
TAIL	1 Bit
VC	2 Bits
TYPE	4 Bits
Payload Data	32 Bytes

Table 3 - Data to Route Subsequent Micropackets in a Message

Description	Size
TAIL	1 Bit
VC	2 Bits
TYPE	4 Bits

Table 4 contains information that can be used to determine whether the micropacket contains errors and a means to report discovered errors.

Table 4 - Data used for error checking and reporting

Description	Size
ERROR	1 Bit
TYPE	4 Bits

Table 4 - Data used for error checking and reporting

		Description	Size
		ECRC	16 Bits
		Payload Data	32 Bytes

Micropacket Transmission order	1	Header micropacket	
	2	1st 32 bytes of user data	
	3	2nd 32 bytes of user data	
		⋮	
	n	Last bytes of user data	

Figure 1 - Message format

Note that there is information used by the switch fabric that also is carried through it.

4.3 Message

As shown in figure 1, Messages are logical groups of micropackets which have the same VC. The first micropacket of a Message, i.e., the Header micropacket, contains information used to route through a HIPPI-6400 fabric (see figure 2) as well as other information as specified in HIPPI-6400-PH. The last micropacket of the Message is marked with the TAIL bit.

5 Switch routing

5.1 Micropacket data transferred through fabric

A HIPPI-6400 switch shall pass the information shown in table 1 through the fabric. Micropacket data payload, the TAIL bit, the TYPE field, the VC field, and the ECRC shall not be modified while passing through a switch fabric. The ERROR bit shall be transferred as set if it was received as set. If the ERROR bit is received as not set, the bit may be set to indicate a switch detected error as described in 5.4.

5.2 Routing of Header micropacket

Figure 2 shows part of the Header micropacket. The complete specification is provided in HIPPI-6400-PH.

Within the Header micropacket, the Destination ULA specifies the Final Destination where a Message is to be sent.

The micropacket TYPE field (TYPE = x'9') identifies a micropacket as a Header micropacket.

TAIL = 1 on a Header micropacket indicates that there are no other micropackets for this Message.

The micropacket VC field specifies one of four logical paths and shall be used to select the appropriate Destination VC (micropackets traverse a fabric on a single VC and never cross VCs).

Switches shall support independent ULA mapping for each input port. This permits mapping the same ULA value to different output ports based upon which input port received the micropacket. See Annex A for an explanation of input port specific switching functionality.

5.2.1 Switch addressing

Switches shall support a mode of operation that provides in-order delivery of all micropackets on a VC from an Originating Source to a Final Destination.

Switches may also provide optional modes of operation such as alternate pathing. These optional modes of operation are not covered by this stan-

Destination ULA DB00-DB05
Source ULA DB06-DB11
Defined in HIPPI-6400-PH DB12-DB31

Figure 2 - Header micropacket addressing

dard and may not guarantee in-order Message delivery.

5.3 Routing of subsequent micropackets in a Message

Subsequent micropackets in a Message (identified by TYPE = x'8' and TYPEs x'B' through x'E') shall be delivered to the same Final Destination that the Header micropacket addressed with its Destination ULA.

The VC field shall be used to distinguish which Message the micropacket belongs to (of the four VCs supported).

When a micropacket is received with the TAIL bit = 1, it indicates that a Message is ended with this micropacket.

5.4 Error protection

If an uncorrectable error is detected in a micropacket that is forwarded, the switch shall set the ERROR bit for that micropacket.

Detected errors shall be logged or counted.

5.4.1 Mandatory error checking

The switch fabric shall pass the unchanged ECRC with each micropacket as specified in HIPPI-6400-PH.

Before sending any micropacket over a HIPPI-6400 link, the switch shall validate the ECRC and set the ERROR bit if the ECRC indicates an error as specified in HIPPI-6400-PH.

5.4.2 Optional error checking

The switch fabric may verify the validity of the ECRC at any point within the fabric.

The switch may also provide additional error detection or correction for internal data errors.

5.4.3 Congestion management

Time-out mechanisms defined in HIPPI-6400-PH will act to prevent switch congestion due to lack of progress on a HIPPI-6400 link, so long as the Source end of the link is functional. However, failures in switch Source ports can prevent this mechanism from functioning.

Switches shall protect against this failure mode by checking Source output ports for continued proper function and by discarding data destined for all failed Source output ports.

5.5 Data interleaving

There are two separate requirements for switch fairness to resolve contention for shared resources. Both micropackets and Messages shall be interleaved as described. These two interleaving processes shall be considered independent and applied without regard to one other.

5.5.1 Micropacket interleaving

Micropacket interleaving between the four VCs shall be applied on a micropacket count basis.

When a switch port has more than one VC with data available for output, the switch shall ensure that micropackets from each VC are afforded an equal opportunity for progress on a physical link.

The algorithm for choosing a micropacket from the available VCs shall allow interleaving on a frequent basis. The recommended algorithm is to interleave VC streams on a single micropacket basis.

Implementations trying to keep short Messages intact (to minimize latency) may use algorithms that interleave on other than a single micropacket basis. No implementations shall permit more than 69 micropackets from a particular VC to be transferred before moving on to the next VC. This limit allows transfer of the maximum permitted VC0 Message (as specified in HIPPI-6400-PH).

Figure 3 shows a simplified switch configuration with two input ports and one output port. Assuming that traffic is available to send to port "C" on more than one VC, a compliant switch alternates between providing output across all busy VCs on link "C", not exceeding the micropacket count limit before switching from one VC to the next VC.

5.5.2 Message interleaving

Message interleaving shall be applied whenever a current Message to an output port is completed.

When a switch has more than one input port with Messages ready for transfer to the same output port (on the same VC), the switch shall ensure that Messages from the input ports are afforded an

equal opportunity for progress. All ports with pending Messages shall be serviced prior to any other port being serviced twice.

In figure 3, an example would be if both port “A” and port “B” have multiple Messages available on their VC0 links ready to send to port “C”. In this example, Messages transferred out VC0 of port “C” are required to alternate between Messages from “A” and “B”.

6 ULA restrictions and reserved ULAs

Although HIPPI-6400 standards provide for a 48-bit ULA space, the total ULA space is not available for all uses. Part of the range of ULAs is reserved to designate the addresses of network services whose location in the network may vary, and for other network management functions.

ULA reservations by the IEEE for network services as described in the Assigned Numbers RFC shall be reserved for the same purposes in HIPPI-6400. Additionally, to allow for backwards compatibility, the 12 bit reserved HIPPI-SC addresses with 36

bits of prefix, as shown in table 5, shall be reserved as ULAs in HIPPI-6400.

All other ULAs are available for assignment to specific Destinations.

These addresses will be changed in the upper bits to show a 48 bit ULA once a block of ULAs has been assigned to the HIPPI Networking Forum.

Note: Later registrations will be added as an addendum to this standard, or as a revision of the standard.

Table 5 - Reserved ULAs

Start of Range	End of Range	Description
x'0F90'	x'0FBF'	Reserved to preserve compatibility with HIPPI-SC address trial self-discovery process
x'0FC0'	x'0FDF'	Reserved for local use

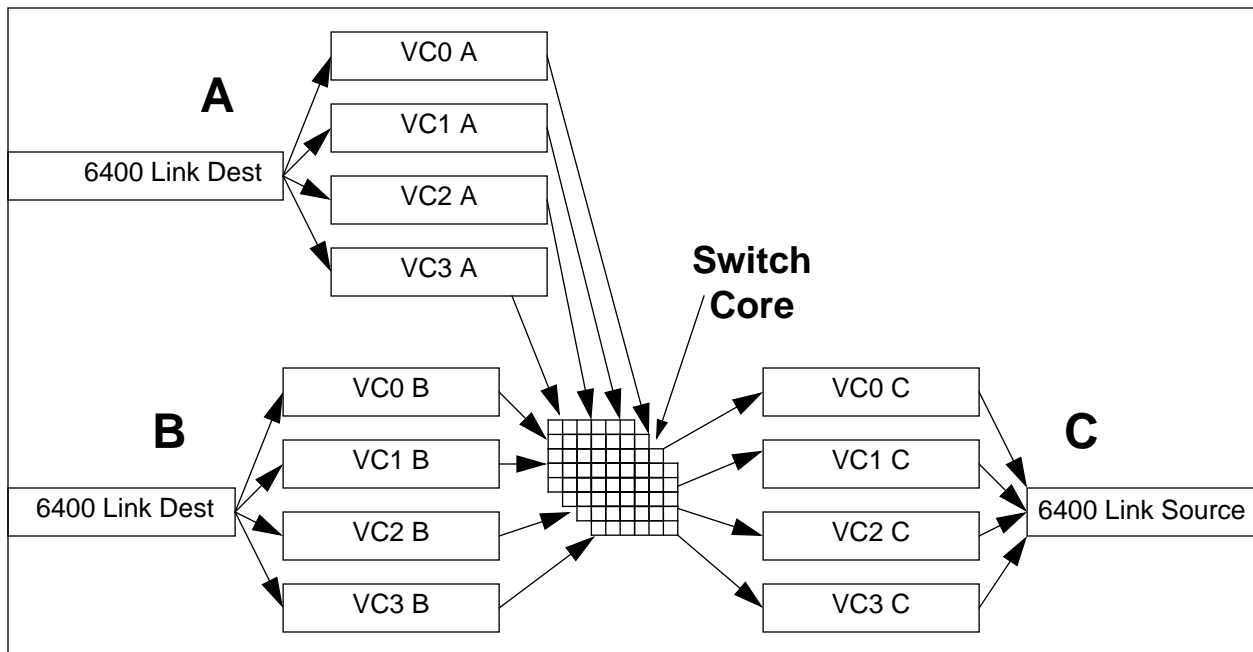


Figure 3 - HIPPI-6400 Switch

Table 5 - Reserved ULAs

Start of Range	End of Range	Description
x'0FE0'	x'0FE0'	Messages pertaining to switch configuration, including HIPPI-LE Address Resolution requests
x'0FE1'	x'0FE1'	All IP protocol traffic conventionally directed to the IEEE 802.1 broadcast address as described in IETF RFC 1042 "Standard for IP transmission over 802.1 networks [2]
x'0FE2'	x'0FE2'	RFC 1112 Host extensions for IP multicasting class D addresses not assigned below [3]
x'0FE3'	x'0FE3'	RFC 1131 OSPF specification All Routers (Class D address 224.0.0.5) [4]
x'0FE5'	x'0FE7'	Reserved
x'0FE8'	x'0FE8'	ISO/IEC 9542:1988 CLNP ES-IS all ES's [5]
x'0FE9'	x'0FE9'	ISO/IEC 9542:1988 CLNP ES-IS all ES's [5]
x'0FEA'	x'0FEA'	ISO/IEC 10589:1992 IS-IS all level 1 IS's [6]
x'0FEB'	x'0FEB'	ISO/IEC 10589:1992 IS-IS all level 2 IS's [6]
x'0FEC'	x'0FEC'	IEEE 802.1d MAC bridging flooding
x'0FED'	x'0FED'	IEEE 802.1d MAC bridging Spanning Tree Protocol
x'0FEE'	x'0FEE'	Embedded switch management agent
x'0FEF'	x'0FFC'	Reserved

Table 5 - Reserved ULAs

Start of Range	End of Range	Description
x'0FFD'	x'0FFD'	Loopback logical address for switches to use when probing other switches
x'0FFE'	x'0FFE'	loopback logical address for hosts to use when probing switches for the host's logical address.
x'0FFF'	x'0FFF'	Unknown or unassigned address. This value should never be used to address a Destination or Destinations. It can be used to indicate that the Source is unaware of its Source address or to signify an unknown logical address in higher layer protocols.

The protocols used to access these services and the means whereby these services keep track of their configuration of the network are outside the scope of this standard.

7 Admin micropackets

Admin micropackets are used for support and initialization of HIPPI-6400 links, Elements, and systems.

I will be adding a diagram to clarify what Elements are at this point in the text.

There are two basic types of Admin micropacket function:

- Within a HIPPI-6400 endpoint or switch, Admin micropackets can be used for internal control of components. This internal usage is done for vendor convenience and is not required to support HIPPI-6400 functionality. Many of the defined Admin micropacket com-

mands will be useful for this control, but the commands used for ULA assignment will not be applicable.

- From one HIPPI-6400 Device (e.g. switch or endpoint) to another, Admin micropackets are used for topology discovery, ULA assignment, and ULA discovery. The ability to send and then receive an echoed micropacket may also be useful as a diagnostic feature. Most other Admin micropacket commands are not useful in this context.

7.1 Element definition

An Element is any component of a HIPPI-6400 system that is able to receive and send micropackets conforming to this standard.

Each end of a HIPPI-6400 link shall operate as an Element. Other components of switches or adapters may optionally conform to the Element definition provided herein. These could include adapter cards, integrated circuits, or software entities.

7.2 Element conformance

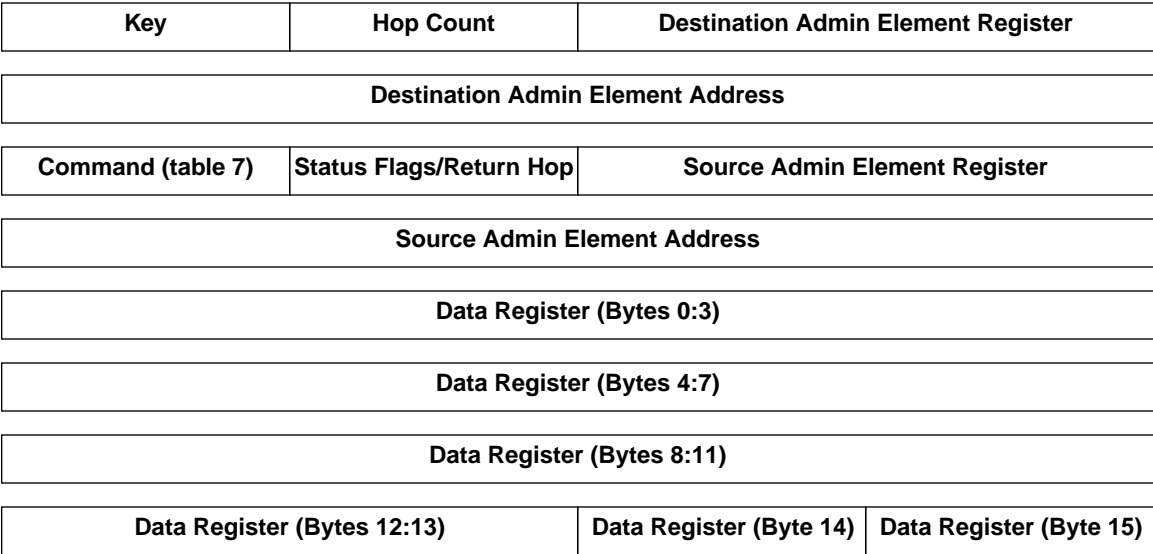
At a minimum, Elements shall support commands and responses for the discovery of Element function, ULA assignment, and ULA discovery. Implementation of other functions called for by Admin

micropacket commands are optional. If an Element does not implement an Admin command, it shall return status to that effect in the response micropacket. All Elements shall respond to each Admin micropacket command with the specified response Admin micropacket.

7.3 Admin micropacket functions

A small set of commands allow for:

- Diagnostic “pings” between HIPPI-6400 Elements, either locally or across a link
- Initial Element address assignment
- Discovery of the function of an Element (e.g. switch or non-switch)
- HIPPI-6400 Source ULA assignment
- Discovery of Destination ULAs attached to a local switch
- Vendor specific register access



Byte 31 of micropacket

Figure 4 - Admin Micropacket Byte Format

7.4 Admin micropacket format

Table 6 and figure 4 both show the format of an

Table 6 - Admin Micropacket Format

Byte	Function
0	Key
1	Hop Count
2:3	Destination Admin Register (designates a local register within an Element)
4:7	Destination Admin Element Address (Destination Element address in a HIPPI-6400 domain)
8	Admin Command (see table 7)
9	Status flags (see table 8) / Return Hop Count
10:12	Source Admin Register (designates a local register within an Element)
12:15	Source Admin Element Address (Source Element address in a HIPPI-6400 domain)
16:31	Data Register

Admin micropacket. Admin micropackets contain:

- **Key:** The Key field is used in certain operations to validate that the originator is authorized to perform the requested operation. Because the key is only 8 bits in length and is returned in response to the SET_ELEMENT_ADDRESS, the protection provided by the key is minimal and only protective against accidental changes. Vendors may also choose to protect their system configuration in other unspecified ways. For example, a vendor may only allow commands that cause configuration changes to occur through a specific port.
- **Hop Count:** If the incoming hop count is zero, the micropacket shall be processed or discarded without a response. If the destination Admin Element address is x'FFFFFFFF', a hop count of zero shall indicate that the Admin micropacket is valid for local processing. All other hop count values in conjunction with a destination Admin Element address of x'FFFFFFFF' indicate that the micropacket shall continue to be forwarded. The value contained in the Hop Count field shall be decremented by

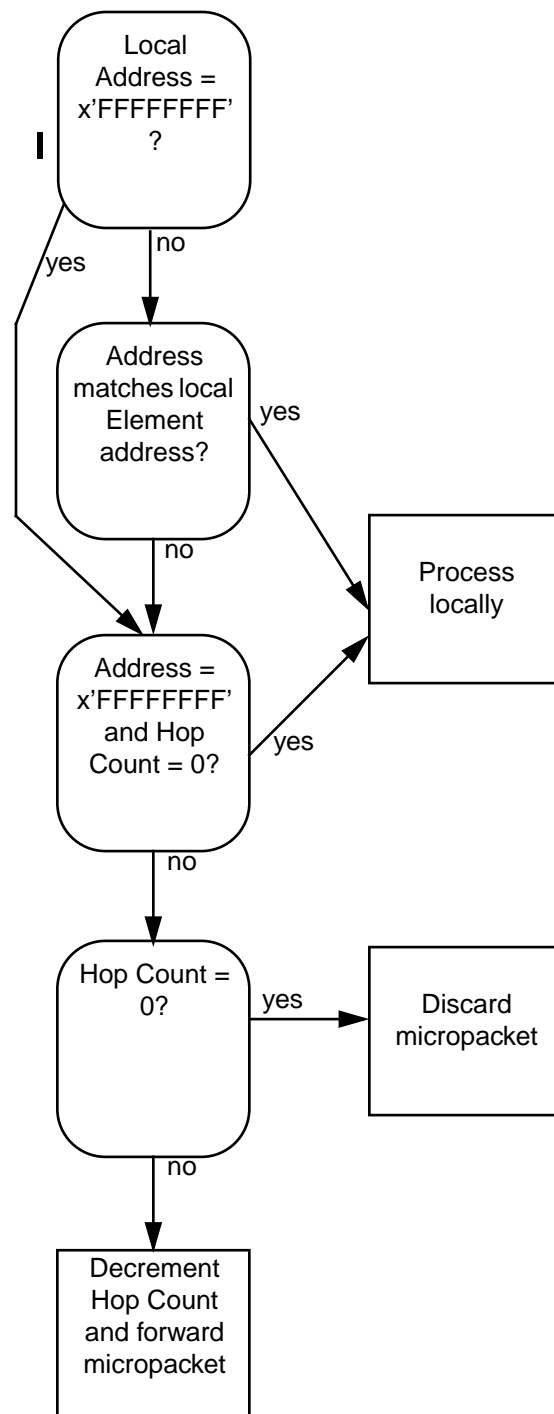


Figure 5 - Admin micropacket

one each time an Admin micropacket exits an Element. If a micropacket is received without a valid Element address match and it cannot be forwarded, it shall be discarded without a response. See figure 5 for a diagram showing how Element addresses are processed.

- Destination Admin Register: The Destination Admin Register field specifies a register within a HIPPI-6400 Element. There are no specific registers required in any Element by this standard and use of any register(s) is optional.
- Destination Admin Element Address: The Destination Admin Element Address field shall be used to specify a particular Element of a HIPPI-6400 system that is the destination of an Admin micropacket command.
- Admin Command: The Admin Command field shall contain a value to specify the meaning and interpretation of the Admin micropacket. Table 7 contains all of the defined values, along with a description of the functions and parameters associated with each command.
- Status Flags / Return Hop Count: When the

Table 8 - Status Flags

Bit	Meaning
0	Undefined Operation
1	Invalid Key
2	Parameter out of range
4	Invalid Element Address
5	Data Register Not valid
6	Unimplemented Command
7	Operation Failed

Admin micropacket is a command, the Return Hop Count field shall be used to communicate the proper hop count value for returning status. The Return Hop Count field may be set to x'FF' when using Element addressing in lieu of a specific return distance.

When the Admin micropacket is a response, the Status Flags field shall be used to return operation results. Table 8 provides definitions for each bit. In each case, flag bit = 1 indicates that the listed exception has occurred.

- Source Admin Register: The Source Admin Register field may be used to specify a register within a HIPPI-6400 Element that can be used as a "reply-to" Element address for certain operations. There are no specific registers required in any Element by this standard.
- Source Admin Element Address: The Source Admin Element Address field is used to specify the particular Element of a HIPPI-6400

system that initiated a sequence of Admin packets. The source Admin Element address shall be used as a "reply-to" Element address.

- Data Register: The Data Register is a 16 byte field that shall be used to carry data for any Admin operation

7.5 Admin micropacket functions

7.5.1 PING

PING may be used to request a response micropacket for diagnostic validation. The Data Register field may be used to send data that will be echoed in the PING_RESPONSE.

The receiving Element shall return a PING_RESPONSE.

7.5.2 PING RESPONSE

PING_RESPONSE acknowledges the PING command. The receiving Element may use this response to validate that the PING'ed Element is operational. The Data Register field shall contain a copy of the data originally sent in the PING command.

7.5.3 SET_ELEMENT_ADDRESS

SET_ELEMENT_ADDRESS may be used to configure an Element with a specific Element address.

The use of Admin micropacket commands for Element address assignment is optional. No Element is required to assign Element addresses.

If this is the first SET_ELEMENT_ADDRESS command received after a reset, the value in the Key field shall be ignored. Later uses of the SET_ELEMENT_ADDRESS command shall validate that the Key field value matches the current key.

If the above criteria for key value are met, the receiving Element shall set its base Admin Element address to be equal to the value set in the lower 4 bytes (12:15) of the Data Register field, shall set its key value to the new key provided in byte 8 of the Data Register. The provided key shall be retained for subsequent command validity checking. Once the base Admin Element address is set, it shall not be changed without validating the key value or until the Element is reset.

Table 7 - Admin Commands

Cmnd Value	Function	V C	Key Req'd	Action	Mandatory?
x'0'	PING	1	No	Asks for a PING_RESPONSE	No
x'1'	PING_RESPONSE	2	No	Acknowledges the PING command	Yes
x'2'	SET_ELEMENT_ADDRESS	1	Yes, except first time after reset	Set base Admin Element address	No
x'3'	SET_ELEMENT_ADDRESS_RESPONSE	2	Yes	Acknowledges the SET_ELEMENT_ADDRESS command	Yes, but may specify that function is not implemented
x'4'	RESET	1	Yes	Commands Element to initialize itself	No
x'5'	EXCHANGE_ELEMENT_FUNCTION	1	No	Provides and requests Element Function	Yes for endpoints and switches, not required for links
x'6'	ELEMENT_FUNCTION_RESPONSE	2	No	Response to a EXCHANGE_ELEMENT_FUNCTION command	Yes
x'7'	ULA_REQUEST	1	No	Requests a Source ULA	Yes for endpoints and switches, no for all others
x'8'	ULA_RESPONSE	2	No	Provides a Source ULA.	Yes for end-points and switches, no for others
x'9'	READ_REGISTER	1	Optional (use of a key may or may not be required)	The sender requests a register value	No
x'A'	READ_REGISTER_RESPONSE	2	No	Returns data from the requested register	No
x'B'	WRITE_REGISTER	1	Optional (use of a key may or may not be required)	Requests that a register value be updated	No
x'C'	WRITE_REGISTER_RESPONSE	2	No	Status for a WRITE_REGISTER	No
x'D'	ERROR_RESPONSE	2	No	Indicates an error	Yes
x'E'	ULA_LIST_REQUEST	1	No	Asks for a list of connected ULAs	No
x'F"	ULA_LIST_RESPONSE	2	No	Provides a list of connected ULAs	Yes for switches
x'10' - x'FF'	Reserved	N / A	N/A	Not defined	No one shall send these

Regardless of the success or failure of the command, the receiving Element shall respond with a SET_ELEMENT_ADDRESS_RESPONSE.

7.5.4 SET_ELEMENT_ADDRESS_RESPONSE

This response acknowledges the SET_ELEMENT_ADDRESS command. The current valid key shall be returned in byte 8 of the Data Register field. The current Element address of this Element shall be returned in the lower 4 bytes (12:15) of the Data Register field. The current Element address and proper key value shall be returned regardless of the success or failure of the SET_ELEMENT_ADDRESS operation.

The use of Admin micropacket commands for Element address assignment is optional. An Element incapable of setting its Element address shall set the Unimplemented Command flag in the flag byte.

7.5.5 RESET

Reset shall cause an Element to initialize itself. This includes clearing the current Element address and key. It may also include other vendor unique functions and may not be the same as the actions caused by a HIPPI-6400 link reset or initialize.

Reset may be propagated further depending upon vendor specific implementation and configuration. There is no message response to a RESET command.

7.5.6 EXCHANGE_ELEMENT_FUNCTION

The sender shall provide its Element function value and set its Element broadcast configuration bits in the least significant byte (0) of the Data Register and requests that the receiver respond with a ELEMENT_FUNCTION_RESPONSE. Element function shall be one of the following:

- Switch Element (00)
Used for a switches that assign logical addresses to endpoints.
- Link-end Element (01)
Used when the Element is a link-end that does not directly assign or use logical addresses
- Non-switch Element (02)
Used for an end-point that requires a logical address assignment.

- Unknown Element (03)

Used when the Element does not deal with logical addresses in any manner, but is an end-point.

The upper two bits in the Element function byte shall be used to signify broadcast parameters for non-switch elements (endpoints). The most significant bit, if set to b'1', shall indicate that this non-switch Element desires to receive broadcast messages. The second most significant bit, if set to b'1', shall indicate that this non-switch Element is capable and willing to act as a broadcast server for this switch.

In order to support switch capability to maintain a list of ports currently desiring broadcast messages, the EXCHANGE_ELEMENT_FUNCTION operation shall be repeated at least once per second and no more than twice per second, when an Element wanting to receive broadcast messages is connected to a responding switch Element (whether broadcast capable or not).

Other bytes in the Data Register are defined as Vendor Unique and may be used in any way desired by the equipment provider.

7.5.7 ELEMENT_FUNCTION_RESPONSE

In response to a EXCHANGE_ELEMENT_FUNCTION command, the sender shall respond by sending the ELEMENT_FUNCTION_RESPONSE with its Element function value in the least significant byte (0) of the Data Register.

Element functions are specified in the EXCHANGE_ELEMENT_FUNCTION command. The most significant bit of the Element function byte shall be echoed as received.

The second most significant bit of the Element function byte shall be set to b'1' if the switch has selected the receiver of this response to be the broadcast server for this switch. The bit will be set to b'0' if the receiver of this response has not been selected to be the broadcast server for this switch.

Other bytes in the Data Register are specified as Vendor Unique and may be used in any way desired by the equipment provider.

7.5.8 ULA_REQUEST

The sender requests that the receiver return a HIPPI-6400 ULA for the sender to use as a Source ULA via a ULA_RESPONSE. The sender shall provide his offered base ULA in bytes (10:15) of the Data Register. The most significant bit of byte 6 of the Data Register shall indicate that this is an additional request for an address and that previous addresses assigned to this port shall be retained as valid in addition to the address(es) assigned by this instance of the command. The balance of bytes (6:7) shall contain a count of desired addresses.

The first ULA registered on a port using the ULA_REQUEST/ULA_RESPONSE is the destination ULA that will be used for broadcast messages.

As specified in 8.2, this command is normally issued by endpoints or switches.

7.5.9 ULA_RESPONSE

The ULA_RESPONSE shall be sent when requested by a ULA_REQUEST command.

Bytes (10:15) of the Data Register shall contain a ULA for the receiver to use as a Source ULA. This may or may not be the offered Source ULA passed in the ULA_REQUEST command. For a receiver needing to add a single Source ULA, this value shall be directly utilized. If the most significant bit of byte 6 is set, it indicates that the Source ULA(s) assigned shall be considered as additional to those assigned in a previous ULA_REQUEST/ULA_RESPONSE operation. For a receiver needing multiple Source ULAs, the balance of bytes (6:7) shall be used as a count of sequential ULAs that start at the base value contained in bytes (10:15) of the Data Register.

Byte (0:5) of the Data Register shall contain a Destination ULA to be used for sending broadcast messages. When a message is sent to this address, it will be broadcast to the entire HIPPI-6400 network.

As specified in 8.2, this response is normally issued by switches.

7.5.10 READ_REGISTER

The sender requests a value from the register specified in the Destination Admin Element Regis-

ter. The receiver shall respond with a READ_REGISTER_RESPONSE.

The use of Admin micropackets for register access is optional. If register access commands are supported, there are no requirements for particular functions or modes specified by this standard.

Contents of registers and their meaning are not specified in this standard.

7.5.11 READ_REGISTER_RESPONSE

The sender shall return the data from the register in the Data Register field.

- Single bytes sent in byte (15)
- Two byte words sent in bytes (14:15)
- Four byte words sent in bytes (12:15)
- Eight byte words sent in bytes (8:15)
- Sixteen byte words sent in bytes (0:15)

The use of Admin micropackets for register access is optional. If register access commands are supported, there are no requirements for particular functions or modes specified by this standard. An Element incapable of supporting this operation shall set the Unimplemented Command flag in the flag byte.

Contents of registers and their meaning are not specified in this standard

7.5.12 WRITE_REGISTER

The sender requests that a register value be updated with the value contained in the Data Register. The receiver shall acknowledge the request with a WRITE_REGISTER_RESPONSE.

- Single bytes sent in byte (15)
- Two byte words sent in bytes (14:15)
- Four byte words sent in bytes (12:15)
- Eight byte words sent in bytes (8:15)
- Sixteen byte words sent in bytes (0:15)

The use of Admin micropackets for register access is optional. If register access commands are supported, there are no requirements for particular functions or modes specified by this standard. No Element is required to issue this command.

Contents of registers and their meaning are not specified in this standard.

7.5.13 WRITE_REGISTER_RESPONSE

The sender shall echo the value written to the specified Data Register. The contents of the Data Register shall be sent as zeros if the update was not successful.

The use of Admin micropackets for register access is optional. If register access commands are supported, there are no requirements for particular functions or modes specified by this standard. An Element incapable of supporting this operation shall set the Unimplemented Command flag in the flag byte.

Contents of registers and their meaning are not specified in this standard.

7.5.14 ERROR RESPONSE

ERROR_RESPONSE shall be sent when an undefined command is received on VC1. No response shall ever be made to Admin micropackets received on VC0, VC2, or VC3.

7.5.15 ULA_LIST_REQUEST

ULA_LIST_REQUEST may be sent to Switch Elements to request a list of attached HIPPI-6400 ULAs. One address may be requested per physical port of the switch Element. This message is sent once for each ULA requested. The data register (byte 2:7) shall contain a number to identify which position in the list is being requested.

7.5.16 ULA_LIST_RESPONSE

ULA_LIST_RESPONSE shall be returned by a switch in response to an ULA_LIST_REQUEST. Switches shall use this response to provide visibility into a sequentially organized list of attached HIPPI ULAs. The list shall contain one entry for each physical port of this switch Element. The first ULA registered using the ULA_REQUEST/ULA_RESPONSE for each port is included in the list. Bytes (10:15) of the Data Register shall contain the ULA. Byte (0) shall contain Function Type for the port. Bytes (2:7) shall contain the list number copied from the ULA_LIST_REQUEST.

In addition to directly connected endpoint ULAs, the sequentially organized list shall also include a single ULA for each directly attached switch. This ULA shall be for the broadcast function of the attached switch.

When access is attempted to list values that have not been validated within the last two seconds (by an EXCHANGE_ELEMENT_FUNCTION operation), the Invalid Element Address bit shall be set in the response. The Operation Failed bit shall not be set in this case.

When access is attempted to list values not included in the list (past the end of the list), the Parameter out of range bit shall be set in the response. The Operation Failed bit shall not be set in this case.

7.5.17 Reserved Admin micropacket functions

Reserved Admin micropacket functions shall not be sent.

Receivers shall perform normal Element address processing and forwarding of Admin micropackets, regardless of the Function code.

Micropackets received for local processing with Reserved Function codes shall be responded to with an ERROR_RESPONSE.

7.6 Addressing of Admin micropackets

The Admin micropacket format contains a 32 bit source and destination Admin Element address. This space is adequate to uniquely identify Elements in configurations of up to 2^{32} Elements.

With Elements that have two ports, a received Admin micropacket shall either be:

- processed locally by the Element
- discarded
- forwarded out the second port

Response Admin micropackets shall be sent on the port that received the original Admin micropacket command. Response Admin micropackets shall use the source Admin Element address and return hop count provided in the original Admin micropacket command as the destination Admin Element address and hop count.

There are two possible destination Admin Element addresses that can result in delivery of an Admin micropacket to an Element for local processing:

- If the destination Admin Element address = x'FFFFFFFF' and hop count = 0. This technique allows access to neighbors (who may possibly have unknown Element addresses) by setting the hop count to control how far distant an Element is in hop count. For example, a hop count of three would pass through three neighboring Elements before being decremented to zero and being processed by the fourth Element.
- When the assigned Element address is not equal to x'FFFFFFFF' and the assigned Element address matches the destination Admin Element address. This technique allows use of a flat logical address space for access to each Element when all of the Element addresses are known.

If a received Admin micropacket contains one of the two possible valid Element addresses pointing to the current local Element, it shall be processed locally. Otherwise, if the hop count value is zero, the packet shall be discarded. Then the hop count shall be decremented by one and the packet shall be forwarded to the Element's other port, i.e., the port that did not deliver this micropacket to this Element.

Admin micropackets shall be sent on the VC specified for each command and response:

- No Admin micropackets shall be sent on VC0 or VC3
- All command Admin micropackets shall be sent on VC1.
- All response Admin micropackets shall be sent on VC2.

Receivers of Admin micropackets shall only process and/or respond to Admin micropackets received on the specified proper VC:

- Admin micropackets received on VC0 or VC3 shall be logged as an error and discarded without a response.
- Admin micropackets received on VC1 shall be processed as a received command, discarded (due to an expired hop count), or forwarded (if the Element address does not match).

- Admin micropackets received on VC2 shall be processed as a received response, discarded (due to an expired hop count), or forwarded (if the Element address does not match). Responses that are received unexpectedly shall be logged as an error and discarded without a response. A response Admin micropacket shall never be sent in reply to an Admin micropacket received on VC2.

Admin micropackets that arrive with either ERROR = 1 or TAIL = 0 shall be logged as an error and discarded without a response.

Selection of the proper port for packet forwarding, from a set of ports in a multi-port Element, is not covered by this standard. Multi-port Element support is optional and may be added in a vendor unique manner.

7.7 Admin Element address assignment

Each Element in a HIPPI-6400 connected collection of Elements may be provided an Element address for operation and control. Element addresses may be assigned through any suitable means, including use of the commands, SET_ELEMENT_ADDRESS and SET_ELEMENT_ADDRESS_RESPONSE. These commands allow an intelligent system Element to assign Element addresses to other Elements within the configuration. Element addresses shall be assigned so that Element address duplication in the connected Element address environment does not occur.

Regardless of whether an Element address is assigned, each Element shall always respond to an Element address of x'FFFFFFFF' when hop count = 0.

This standard does not specify how the intelligent system Element chooses Element addresses for assignment. The discovery of topologies beyond two ports and the mechanisms for multi-port Element address assignment are not covered by this standard. Multi-port Element support is optional and may be added in a vendor unique manner.

7.8 Admin micropacket flow control

Admin micropacket operations (with the exception of reset) consist of a command and a paired response operation. To avoid overrun of receivers, no more than one operation shall be outstanding to

a single destination Element from a single source Element in a time period of one second. Therefore, Elements shall send only a single command:

- PING
- SET_ELEMENT_ADDRESS
- EXCHANGE_ELEMENT_FUNCTION
- ULA_REQUEST
- READ_REGISTER
- WRITE_REGISTER
- ULA_LIST_REQUEST

before receiving the paired response micropacket:

- PING_RESPONSE
- SET_ELEMENT_ADDRESS_RESPONSE
- ELEMENT_FUNCTION_RESPONSE
- ULA_RESPONSE
- READ_REGISTER_RESPONSE
- WRITE_REGISTER_RESPONSE
- ULA_LIST_RESPONSE

or until a time-out period of at least one second has elapsed.

Since RESET has no response, Elements that have sent a RESET shall wait at least one second before attempting any other operation to the Element that has been reset.

8 ULA Configuration

In addition to switching HIPPI-6400 Messages between ports, HIPPI-6400 ports shall support in-band communications for switch management functions.

To support topology discovery and ULA configuration, HIPPI-6400 Destination ports shall be capable of receiving and processing micropackets with TYPE = Admin over any connected HIPPI-6400 link.

To support topology discovery and ULA configuration, HIPPI-6400 Source ports shall be capable of

sending micropackets of TYPE = Admin over any connected HIPPI-6400 link.

8.1 Determination of Topology

As a step in the procedure to establish a ULA for self identification (used as the Source ULA field), endpoints and switches shall identify if they are connected to another endpoint or to a switch.

Intervening link support hardware and interface components may be present on either side of a HIPPI-6400 link. These intermediate Elements will typically not contain information useful for ULA assignment. The endpoint discovering topology information shall identify these intermediate points to discover the location of an Element capable of exchanging information about ULA configuration.

Information about the function of connected Elements is collected by sending an EXCHANGE_ELEMENT_FUNCTION Admin micropacket. The endpoint may directly select a destination if the appropriate Admin Element address information is already known, or it may use hop-count Element addressing to discover what is connected and how far away (in hops) the Element of interest is located.

If an Element responds that it is a link Element or an unknown Element, the probing system shall continue to the next Element. Once a connected Element is identified as an endpoint or switch, topology determination is complete.

In figure 6, an example of an endpoint to endpoint link is shown. In this example, System A needs to determine the Element function of System B, for ULA configuration. System B also needs to determine the Element function of System A, for the same reason. The following example traces the operation of System A.

System A begins by probing each Element that supports Admin micropackets until it reaches the endpoint of System B.

- a. System A sends an EXCHANGE_ELEMENT_FUNCTION Admin micropacket to the closest point with an Element address of x'FFFFFFF' and a hop-count of 0. This will be received and processed by Link-End A. Link-End A will respond in the ELEMENT_FUNCTION_RESPONSE Admin

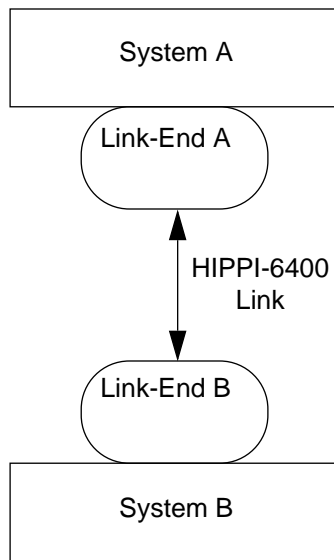


Figure 6 - Endpoint to endpoint connect

micropacket that it is a link. System A must therefore go further to reach another endpoint or switch.

- b. System A sends an `EXCHANGE_ELEMENT_FUNCTION` Admin micropacket to the next closest point with an Element address of `x'FFFFFFFF'` and a hop-count of 1. This will be received and processed by Link-End B. Link-End B will respond in the `ELEMENT_FUNCTION_RESPONSE` Admin micropacket that it is a link. System A must therefore go further to reach another endpoint or switch.
- c. System A sends an `EXCHANGE_ELEMENT_FUNCTION` Admin micropacket to the 3rd closest point with an Element address of `x'FFFFFFFF'` and a hop-count of 2. This will be received and processed by System B. System B will respond in the `ELEMENT_FUNCTION_RESPONSE` Admin micropacket that it is an endpoint. System A now knows where to exchange information regarding ULAs.

In the above example, System A presumably would have been aware that the most directly attached component (Link-End A) is part of its own configuration and that it need not communicate with that component. It therefore would not have needed to start with a hop-count of 0 (but was not detrimentally effected by doing so).

System B could determine the type of Element for System A in two ways:

- System B could duplicate the above steps in reverse.
- System B could use the information provided in the `EXCHANGE_ELEMENT_FUNCTION` command that System A sent to System B when the third step in the above exchange took place. Endpoints should not wait for the other end to perform an exchange, but if the exchange occurs at an appropriate time, they may take advantage of the occurrence.

8.2 ULA exchange

Once the other end of the link has been identified as to type (switch or non-switch endpoint), ULAs are configured.

8.2.1 Endpoints on both ends

If both ends of the link are endpoints, each side shall use any desired value(s) for a Source ULA. Although duplicate Source ULAs between the two ends of the link are possible, having the same ULA(s) at both ends will not prevent proper HIPPI-6400 operation.

If the other end of a link cannot be determined (the last reachable Element does not respond properly as a switch or endpoint), the non-responsive end shall be treated as an endpoint.

8.2.2 Switches on both ends

Switch to switch ULA configuration shall occur to exchange ULAs for the broadcast function. A `ULA_REQUEST` shall be sent by each switch to all directly connected switches. Upon receipt of a `ULA_REQUEST` Admin micropacket, the receiving switch shall respond with a `ULA_RESPONSE` Admin micropacket. The `ULA_RESPONSE` shall contain a ULA that is valid either as a ULA for an embedded broadcast capability or is routed to a broadcast server.

Switch to switch ULA discovery to learn the full set of connected ULAs on distant switches is handled outside of this standard. Methods of switch configuration could include static manual table entry or automated ULA learning algorithms.

8.2.3 Endpoint to switch

If endpoints discover that they are connected to switches, they shall advertise a Source ULA. The ULA offer shall be made by sending a ULA_REQUEST Admin micropacket.

Mechanisms for selection of this advertised ULA are not specified by this document and any desired approach may be used. A common approach of network equipment vendors is to use a ULA from a block of ULAs purchased from the IEEE. This method provides some guarantee of uniqueness and is recommended unless there are factors that require a different approach.

Upon receipt of a ULA_REQUEST Admin micropacket, the receiver shall respond with a ULA_RESPONSE Admin micropacket. The ULA_RESPONSE shall contain a Source ULA valid for the ULA_RESPONSE recipient. This ULA may be the same as advertised in the original ULA_REQUEST offer or it may be different.

This returned Source ULA shall be accepted and subsequently used in all HIPPI-6400 messages by the receiver of the ULA_RESPONSE Admin micropacket.

Regardless of whether the returned Source ULA is the same as the Source ULA originally offered by the endpoint, the switch is the final selector of the Source ULA that will be used by the endpoint.

Switches shall prevent a single ULA from being assigned more than once in the same fabric.

9 Broadcast

All switch Elements shall either directly support broadcast of messages or shall provide support of broadcast servers.

Any host that requires broadcast functionality should implement a broadcast server function to guarantee that broadcast functionality will be available if connected to a non-broadcast capable switch.

9.1 Broadcast Operation

Each switch shall provide a unique broadcast Destination ULA. This ULA is advertised to any attached port that performs a registration of a Source ULA.

Messages sent to the broadcast address shall be delivered to all hosts within a HIPPI-6400 fabric that have registered their desire to receive broadcasts.

Broadcast capable switches shall deliver broadcast messages directly to each attached port.

Non-broadcast capable switches shall route broadcast messages to a broadcast server. Endpoints selected as broadcast servers shall forward received broadcast messages, so that broadcast messages are delivered to each attached port that have registered to receive them.

9.2 Selection of broadcast server

Non-broadcast capable switches shall select a broadcast server from attached hosts who have indicated their capability and willingness to perform the broadcast server function. The indication of qualified broadcast servers is provided by an EXCHANGE_ELEMENT_FUNCTION operation with the second most significant bit of the Element function byte set to b'1'.

One server shall be selected per non-broadcast capable switch. The server shall be notified by returning the second most significant bit of the Element function byte set to b'1' in the ELEMENT_FUNCTION_RESPONSE.

The ELEMENT_FUNCTION_EXCHANGE and ELEMENT_FUNCTION_RESPONSE shall be exchanged at intervals of from one to two per second. This continued exchange allows selection of a broadcast server as needed to deal with equipment failures and to accommodate added or removed systems. If a broadcast server fails to provide a ELEMENT_FUNCTION_EXCHANGE message within 5 seconds, the switch shall select a new broadcast server.

9.3 Registration for broadcast

9.3.1 Endpoints

Attached endpoints may register to receive broadcasts. This shall be done by setting the most significant bit of the Element function byte to b'1' in the ELEMENT_FUNCTION_EXCHANGE operation. To maintain registration for receiving broadcast messages, the ELEMENT_FUNCTION_EXCHANGE shall be repeated at a rate of one to two times per second.

Attached endpoints may choose not to receive broadcast messages. This shall be done by either:

- not registering for broadcast messages
- re-registering with the ELEMENT_FUNCTION_EXCHANGE operation and setting the most significant bit of the Element function byte to b'0'
- allowing the broadcast registration to time-out

9.3.2 Switches

Switches are required to maintain a list of ULAs for broadcast. This list shall include one entry for:

- each endpoint directly connected to this switch that has made at least one ULA_REQUEST and has registered a request to receive broadcast messages within the last 5 seconds through the ELEMENT_FUNCTION_EXCHANGE process
- each other switch directly connected to this switch that has provided a broadcast address via the ULA_REQUEST/ULA_RESPONSE process

The list of broadcast destinations shall be made available through the ULA_LIST_REQUEST and ULA_LIST_RESPONSE.

More descriptions or pointers need to be added on how broadcast works across switches.

Annex A (informative)

Switching

A.1 General

HIPPI-6400 switching of Messages is accomplished by processing the Destination ULA field of the HIPPI-6400-PH MAC header. This may be done based on the complete contents of the Destination ULA (48 bits) or on a subset of the field.

If a subset of the Destination ULA is used for switching, switches must ensure that Source ULAs are unique in the portion of the ULA operated on by the switch. Clause 8 describes the process of ULA configuration that gives switches final authority in configuration of Source ULAs.

When connections are made to other networks, the address range of the two (or more) networks is limited by the smaller of the connected address ranges.

For example, HIPPI-PH can be switched to communicate with HIPPI-6400 so long as all of the communicating systems restrict their addresses to 12 bits. The total number of addresses is therefore limited to 4096 (minus reserved addresses).

The Destination ULA field in the Header micro-packet is used to control HIPPI-6400 physical layer switches, supporting the interconnection of many Devices. Figure 7 shows an example configuration that will be used to describe how HIPPI-6400 switches function. Three hosts and two switches are shown, actual configurations may be smaller or larger.

Although there is only a single mode of operation (ULA addressing) specified for HIPPI-6400, users can achieve a form of source routing (as described in HIPPI-SC) by their selection of port configuration.

A.2 Logical addressing

With logical addressing, ULAs specify where a Message is to be delivered, not the path to take to get there. Originating Sources use the same ULA to reach a Final Destination, no matter where the Originating Source is located.

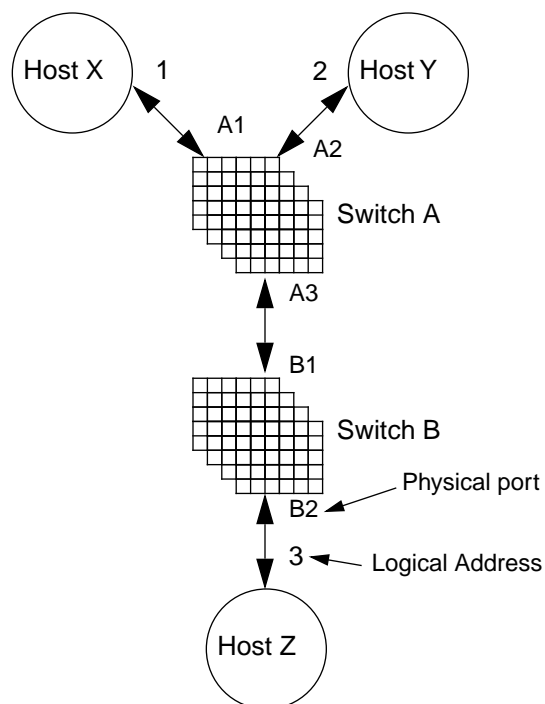


Figure 7 - Hosts and switch configuration

In figure 7, Host X, Host Y, and even Host Z can use ULA "3" to specify that a Message should be sent to Host Z.

With ULAs, the intermediate switches are responsible for selecting an appropriate path.

It is envisioned that switches can be built to use look-up tables at each input port to map ULAs to Destinations. A look-up table can be indexed using the Destination ULA field. The look-up table would be used to hold a possible path(s) for a Destination.

A major advantage of using ULAs is that only the switches need to know the fabric interconnection topology and the hosts only need to know the ULAs. Hence if a link or port fails, switches can address around it without the hosts having to know about it or do anything special.

A.3 Input specific logical addressing

Because each input port is specified to contain a unique ULA look-up capability, it is possible to use logical switch addressing for limited source routing. Note that only the input portion of a port is involved in addressing. When a Message exits on a particular output port, it crosses that link without further addressing until received at the next input.

This capability means that it is possible to create addressing that could result in infinite looping of a micropacket. This will rarely be desirable and should be avoided.

One possible use of input port specific routing is to provide a test capability for monitoring the performance of specific links. In figure 7, if Host Y wants to monitor the state of the link between switch A and switch B, he can send a Message to switch A

and then to switch B. Port B1's ULA table (at switch B) can direct the Message back to B1, then switch A, and back to Host Y. To do this, the same ULA must be handled differently by individual ports. Table 9 shows a simplified look-up table that would work in this example.

Table 9 - Port look-up table

ULA	Port Number	Destination
2	A2	A3
2	B1	B1
2	A3	A2

Because there are many available ULAs, normal flat addressing can be used for host communications with other ULAs used to support input specific logical routing for test and monitoring purposes.

Annex B (informative)

Bridging

B.1 General

I believe that bridging with HIPPI-6400, as described here is no longer possible (the destination ULA is changed). Is this true? Is there something else that needs to go here?

HIPPI-6400 bridging may be used as a substitute for directly manipulating MAC ULAs of connected media types. Bridges use ULAs embedded in the Message body as a look-up for the current media address.

With bridging, the incoming ULA is only used to send across a single fabric. At each media translation to and/or from HIPPI-6400, a new MAC address is found based on the Message ULA address.

For example, in figure 8, Host X uses a ULA of “4” to communicate with the bridge system when sending a message to Host W (ULA “6” on a separate network). The bridge operates on a ULA contained in the message body to look up the address of Host W on the connected network. Host W would reply using the bridge ULA (“5”) and the bridge would look up the ULA of X (“1”) to place in the HIPPI-6400 MAC header

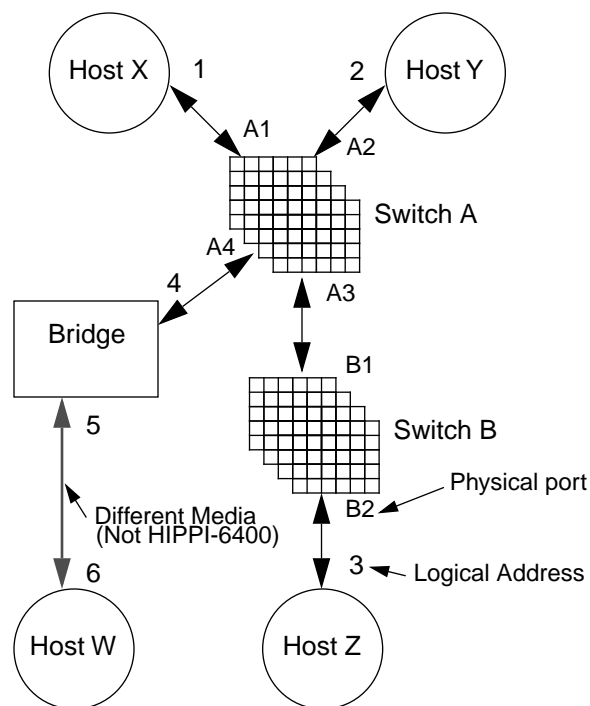


Figure 8 - Hosts, switch, and bridge configuration

Building of look-up tables for bridging operations can be done using an automated process such as ARP or can be handled with static table entries.

Annex C (informative)

Bibliography

The following documents are the basis for assignment of specific logical addresses for certain network services.

- [1] RFC 1042, Standard for the transmission of IP datagrams over IEEE 802 networks. (Provides the general techniques that the Internet Protocol uses to build media packet headers on IEEE 802 (IS 8802) networks.)
- [2] RFC 2067, IP on HIPPI. (Describes a technique whereby hosts may use the Internet Protocol over a HIPPI compliant interface.)
- [3] RFC 1112, Host extensions for IP multicasting. (Provides a technique whereby network and transport layer Internet protocol applications may use the multicasting capabilities defined for IS 8802 networks.)
- [4] RFC 1131, OSPF specification. (Describes the open shortest path first IP protocol which permits network layer IP routers to discover the best route to remote IP addresses and networks.)
- [5] ISO/IEC 9542:1988, Telecommunications and information exchange between systems – End system to intermediate system routing exchange protocol for use in conjunction with the protocol for providing the connectionless-mode network service (ISO 8473)
- [6] ISO/IEC 10589:1992, Telecommunications and information exchange between systems – Intermediate system to intermediate system intra-domain routing exchange protocol for use in conjunction with the protocol for providing the connectionless-mode network service (ISO 8473)
- [7] ANSI/IEEE 802.1D-1990, Media access control (MAC) bridges, (Specifies the operation of transparent bridges between IEEE 802 conformant networks.)

Note: RFC (Request For Comment) documents are working standards documents from the TCP/IP internetworking community. Copies of these documents are available from numerous electronic sources (e.g., <http://www.ietf.org>) or by writing to IETF Secretariat, c/o Corporation for National Research Initiatives, 1895 Preston White Drive, Suite 100 Reston, VA 20191-5434, USA.